

Parallel Implementation of Interval Matrix Multiplication *

Nathalie Revol

INRIA Grenoble - Rhône-Alpes – Université de Lyon
LIP (UMR 5668 CNRS - ENS Lyon - UCB Lyon 1 -
INRIA), ENS de Lyon, France
`nathalie.revol@ens-lyon.fr`

Philippe Théveny

École Normale Supérieure de Lyon–Université de Lyon
LIP (UMR 5668 CNRS - ENS Lyon - UCB Lyon 1 -
INRIA), ENS de Lyon, France
`philippe.theveny@ens-lyon.fr`

Abstract

Two main and not necessarily compatible objectives when implementing the product of two dense matrices with interval coefficients are accuracy and efficiency. In this work, we focus on an implementation on multicore architectures. One direction successfully explored to gain performance in execution time is the representation of intervals by their midpoints and radii rather than the classical representation by endpoints. Computing with the midpoint-radius representation enables the use of optimized floating-point BLAS and, consequently, the performances benefit from the performances of the BLAS routines. Several variants of interval matrix multiplication have been proposed, which correspond to various trade-offs between accuracy and efficiency, including some efficient ones proposed by Rump in 2012. However, in order to guarantee that the computed result encloses the exact one, these efficient algorithms rely on an assumption on the order of execution of floating-point operations, which is not verified by most implementations of BLAS. In this paper, an algorithm for interval matrix product is proposed that verifies this assumption. Furthermore, several optimizations are proposed and the implementation on a multicore architecture compares reasonably well with a non-guaranteed implementation based on MKL, the optimized BLAS of Intel: the overhead is less than 2 for matrix size up to 3,500. This implementation also exhibits a good scalability.

Keywords: interval arithmetic, interval matrix multiplication, BLAS, gemm, parallel computer, multicore

AMS subject classifications: 65Y05, 65G40, 65F30, 65-04

*This work is partly funded by the HPAC project of the French Agence Nationale de la Recherche (ANR 11 BS02 013).

1 Introduction

When implementing the multiplication of two dense matrices with interval coefficients, two conflicting objectives must be met. A first objective is to get accurate results, i.e. results that do not overestimate too much the exact results, and preferably with known and bounded overestimation ratios. A second objective is to get results quickly, in this work particularly, we want to exploit the capabilities of multicore architectures and of the corresponding programming paradigm, namely multithreading. Several algorithms, that offer different trade-offs between these two objectives, are established by Nguyen in [12].

The difficulties one has to face when implementing an interval matrix multiplication are manifold. Implementing interval arithmetic through floating-point arithmetic relies on changes of the rounding modes, either rounding downwards and upwards with the representation by endpoints, or rounding to-nearest and upwards with the so-called midpoint-radius representation, using the midpoints and radii. Whether the rounding mode is kept unchanged or modified by BLAS routines is undocumented. Furthermore, whether the rounding mode is properly saved and restored at context switches, as in a multithreaded execution, is not documented either, as pointed out in [8]. Another issue in numerical, floating-point computing, which also impacts the behavior of routines in interval arithmetic, is the non-reproducibility of computations, especially on HPC systems, as noted in [4]. This phenomenon originates, in particular, from the fact that the floating-point addition is not associative and thus that the order of the additions performed in, for instance, the sum of many numbers depends on the number of threads, the scheduler, the storage in memory of these numbers etc. Thus it is not realistic to rely on any assumption on the order of floating-point computations, which is a problem for the applicability of Theorem 2.1 in Section 2. A last kind of difficulties is related to the programming of multicore architectures. As optimized libraries, such as MKL [6], optimized by Intel for its processors, exhibit good performances, a first idea is to try to use these libraries to benefit from their performances. However, as these libraries do not fulfill any assumption on the execution order of their operations, as mentioned above, they cannot be used and one turns to coding and optimizing the code by hand, using for instance OpenMP and “helping” the compiler to vectorize the code. Getting good performances in the development of such codes is difficult.

In this paper, we focus on algorithms for interval matrix multiplication. Various formulas are based on the midpoint-radius representation, ranging from exact ones, as in [11, pp. 22-23], to faster ones in [15] that resort to 4 calls to BLAS floating-point matrix multiplication, including intermediate ones in terms of accuracy and number of calls to BLAS routines given in [12, Chapter 2]. Algorithms that use interval matrix multiplication to verify a floating-point matrix product are given in [14], algorithms that save some floating-point matrix products by bounding the roundoff errors are given in [16]. We will focus on the latter ones. BLAS routines for the product of dense floating-point matrices are detailed in [2]. Optimized versions, such as Goto’s BLAS [3] or Intel’s MKL, yield very good performances, autotuned versions, such as ATLAS [19], provide good performances on various hardware support by tuning automatically cache usage via block size etc.: they are good candidates for interval matrix multiplication algorithms, such as those mentioned above, that resort to calls to BLAS floating-point matrix multiplication.

We also focus on architectures that become more and more frequent nowadays and very probably even more in the near future, namely multicore architectures. The machine we use in our experiments is a multiprocessor where each processor is composed of several cores, for even higher performances. This architecture still relies on a shared memory: going to distributed architecture is a major step we did not take, as issues become fundamentally different from shared-memory or sequential ones.

We recall in Section 2 the algorithms from [16] and the assumption on which they rely. As this assumption is no more verified by the implementation of BLAS on the architecture we consider, as explained in Section 3, we give up the use of BLAS. The main result of this paper is **the implementation of dense interval matrix multiplication which is fully guaranteed** as the requirement is fulfilled, **and which gives good performances** as hand-made optimizations are done where needed. It is detailed in Section 4. In particular, exploiting the cache usage through a version by blocks yields performances which are close to the performance of a (non-guaranteed) implementation based on Intel's MKL, in terms of execution time and scalability, as shown in Section 5.

2 Algorithms for Interval Matrix Multiplication

In this paper we use the notation for real interval analysis defined in [7]. A point matrix has its coefficients which are punctual: real or floating-point numbers, as opposed to an interval matrix whose coefficients are intervals. An interval matrix $\mathbf{A} \in \mathbb{IR}^{n \times m}$ can be written as a pair of point matrices $[\underline{\mathbf{A}}, \overline{\mathbf{A}}]$ and represents the set of all matrices $M \in \mathbb{R}^{m \times n}$ such that $\underline{\mathbf{A}} \leq M \leq \overline{\mathbf{A}}$, where inequalities hold componentwise. Equivalently, the matrix \mathbf{A} can be represented by the pair of point matrices $\langle \text{mid } \mathbf{A}, \text{rad } \mathbf{A} \rangle$ where the (i, j) component of the midpoint matrix $\text{mid } \mathbf{A}$ (respectively, of the radius matrix $\text{rad } \mathbf{A}$) is the midpoint (respectively, the radius) of $[\underline{A}_{ij}, \overline{A}_{ij}]$:

$$\text{mid } \mathbf{A}_{ij} = \frac{\underline{A}_{ij} + \overline{A}_{ij}}{2} \quad \text{and} \quad \text{rad } \mathbf{A}_{ij} = \frac{\overline{A}_{ij} - \underline{A}_{ij}}{2}.$$

In [15], Rump showed that the midpoint-radius representation of matrices enables faster matrix operations on computer than their infimum-supremum counterparts. In particular, he described an interval matrix product algorithm that relies on point matrix multiplications for its most expensive computational part. The gain in execution time with this approach comes from the fact that possibly costly rounding mode changes are limited and that numerous implementations of BLAS provide efficient point matrix products. This midpoint-radius algorithm requires four point matrix products and is proven to compute an enclosure of the exact matrix product with a radius that is overestimated by a factor¹ at most 1.5.

In [13], Ogita and Oishi proposed faster midpoint-radius algorithms computing the interval matrix product in approximately twice the cost of a point matrix product. However, the known bound for the radius overestimation is worse than with the previous algorithm.

¹Roundoff errors are neglected in radius overestimation factors given here. This is acceptable when the input intervals are large enough, say $\text{rad } \mathbf{x} > 10^{-12} \text{ mid } \mathbf{x}$ for double precision.

In his PhD thesis [12, Chapter 2], Nguyen exhibited another interval matrix product algorithm that overestimates the radius by a factor less than 1.18 at a computational cost less than twice the initial midpoint-radius algorithm from Rump.

In [16], Rump improved the control of rounding errors in the product of midpoint matrices with the following bound (the inequality applies componentwise):

Theorem 2.1 *Let $A \in \mathbb{F}^{m \times k}$ and $B \in \mathbb{F}^{k \times n}$ with $2(k+2)\mathbf{u} \leq 1$ be given, and let $C = fl_{\square}(A \cdot B)$ and $\Gamma = fl_{\square}(|A| \cdot |B|)$. Here, C may be computed in any order, and we assume that Γ is computed in the same order. Then*

$$|fl_{\square}(A \cdot B) - A \cdot B| \leq fl_{\square} \left(\frac{k+2}{2} \text{ulp}(\Gamma) + \frac{1}{2} \mathbf{u}^{-1} \eta \right).$$

In this theorem and in what follows, \mathbb{F} denotes the set of floating-point numbers, \mathbb{IF} the set of intervals with floating-point (midpoint-radius in our case) representation, \mathbf{u} and η respectively the unit roundoff and the smallest positive normal floating-point number, $\text{ulp}(\Gamma)$ the unit in the last place² of the components of Γ , and $fl_{\square}(E)$ (resp. $fl_{\Delta}(E)$) indicates that every operation in the arithmetic expression E is evaluated with rounding-to-nearest (resp. rounding towards $+\infty$) mode. Using this bound, Rump transformed the above algorithms eliminating one point matrix product from his previous interval matrix product and two from Nguyen's. We reproduce here the improved algorithms **MMu13** (Algorithm 1) and **MMu15** (Algorithm³ 2) without the additional conversions from and to the usual infimum-supremum representation of matrices.

Algorithm 1 **MMu13**

Input: $A = \langle M_A, R_A \rangle \in \mathbb{IF}^{m \times k}, B = \langle M_B, R_B \rangle \in \mathbb{IF}^{k \times n}$

Output: $C_3 \supseteq A \cdot B$

- 1: $M_C \leftarrow fl_{\square}(M_A \cdot M_B)$
 - 2: $R'_B \leftarrow fl_{\Delta}((k+2)\mathbf{u}|M_B| + R_B)$
 - 3: $R_C \leftarrow fl_{\Delta}(|M_A| \cdot R'_B + R_A \cdot (|M_B| + R_B) + \mathbf{u}^{-1}\eta)$
 - 4: **return** $\langle M_C, R_C \rangle$
-

Algorithm 2 **MMu15**

Input: $A = \langle M_A, R_A \rangle \in \mathbb{IF}^{m \times k}, B = \langle M_B, R_B \rangle \in \mathbb{IF}^{k \times n}$

Output: $C_5 \supseteq A \cdot B$

- 1: $\rho_A \leftarrow \text{sign}(M_A) \cdot \min(|M_A|, R_A)$
 - 2: $\rho_B \leftarrow \text{sign}(M_B) \cdot \min(|M_B|, R_B)$
 - 3: $M_C \leftarrow fl_{\square}(M_A \cdot M_B + \rho_A \cdot \rho_B)$
 - 4: $\Gamma \leftarrow fl_{\square}(|M_A| \cdot |M_B| + |\rho_A| \cdot |\rho_B|)$
 - 5: $\gamma \leftarrow fl_{\Delta}((k+1)\text{ulp}(\Gamma) + \frac{1}{2}\mathbf{u}^{-1}\eta)$
 - 6: $R_C \leftarrow fl_{\Delta}((|M_A| + R_A) \cdot (|M_B| + R_B) - \Gamma + 2\gamma)$
 - 7: **return** $\langle M_C, R_C \rangle$
-

²see [10, Section 2.6] for a definition.

³Operations \cdot in lines 1 and 2 are the componentwise multiplication, as in Matlab notation.

Finally, Ozaki et al. in [14] used coarser but faster estimates for the radius of the product and proposed several algorithms for interval matrix multiplication at the cost of one, two and three calls to point matrix multiplication.

By relying on floating-point matrix multiplications, the previous algorithms were designed to be easy to implement efficiently. However, some care is needed and the next section discusses some implementation issues that can easily be overlooked.

3 Implementation Issues

The midpoint-radius representation of matrices allows one to resort to floating-point matrix products for the implementation of interval matrix product. Many software libraries (for example GotoBLAS [3] or ATLAS [19]) offer optimized implementations of matrix operations following the interface defined by the Basic Linear Algebra Subprograms Technical Forum Standard [2]. Processor vendors also provide BLAS libraries with state-of-the-art implementation of matrix products for their own architectures, like, for instance, the Intel Math Kernel Library (MKL) and the AMD Core Math Library. So, BLAS libraries give access to high-performance and portability.

However, `MMul3` (Algorithm 1) and `MMul5` (Algorithm 2) require that two hypotheses are satisfied: first, the rounding mode is taken into account by the library and, second, $M_A \cdot M_B$ and $|M_A| \cdot |M_B|$ are computed in the same order (along with $\rho_A \cdot \rho_B$ and $|\rho_A| \cdot |\rho_B|$ in `MMul5`).

3.1 Rounding Modes

We mention here some situations where implementations with BLAS fall short of the first hypothesis when the rounding mode in use is not the default rounding-to-nearest (see [8] for a thorough discussion). For instance, in order to calculate an overestimated result, a matrix product using a Strassen-like recursive algorithm [17] should compute overestimations for terms that are added and underestimations for terms that are subtracted, but this would ruin its advantage over the classical algorithm. Extended precision BLAS [2, Chapter 4] are another example. For matrix operations especially, the reference implementation [9] uses double-double arithmetic which makes intensive use of error-free transformations [10, Chapter 4] that require rounding-to-nearest mode.

Nonetheless, the radius of the interval product in Algorithms 1 and 2 above must be computed with rounding towards $+\infty$ to guarantee that the result accounts for all possible roundoff errors and contains the exact product.

3.2 Computation Order of $M_A \cdot M_B$ and $|M_A| \cdot |M_B|$

The second hypothesis comes from Theorem 2.1. The BLAS interface does not provide any function that calculates simultaneously and in the same order a matrix product and the product of the absolute values of the inputs. So, the implementations of midpoint-radius algorithms have to compute the quantities of Theorem 2.1 with successive calls to floating-point matrix products, as in lines 3 and 4 of Algorithm 2. However, no order of operation is specified by the BLAS standard [2], so Theorem 2.1 may not be applicable. The orders of the floating-point operations performed by the two multiplications are even more likely to differ when the BLAS library is itself multithreaded.

In that respect, it is useful to note that legal obligations and verification constraints of their users cause some vendors to address the problem of the reproducibility of numerical results between different processors or from run to run. For these reasons, the version 11.0 of MKL (see [18] or [6, Chapter 7]) now provides modes of execution where the user can control, at some loss in efficiency, the task scheduling and the type of computing kernels in use. In these modes, identical numerical results are guaranteed on different processors when they share the same architecture and run the same operating system. Moreover, reproducibility from run to run is ensured under the condition that, in all executions, the matrices have the same memory alignment and the number of threads remains constant.

We can use this kind of control to solve the problem of the computation order of $M_{\mathbf{A}} \cdot M_{\mathbf{B}}$ and $|M_{\mathbf{A}}| \cdot |M_{\mathbf{B}}|$. As the processor and the OS remain the same during the computation of the two products computations, it suffices, first, to compute $M_{\mathbf{A}} \cdot M_{\mathbf{B}}$, second, to transform in place matrix components into their absolute values ensuring the identity of memory alignments, then third, to recompute the product on the new input with the same number of threads. The drawback of this solution is its specificity to the Intel MKL, as long as other libraries do not adopt the same control for numerical reproducibility.

In order to overcome these uncertainties on the execution of an arbitrary BLAS library, we present in the next section several implementations of Algorithm **MMMu15** that verify the two assumptions while still being efficient on a multicore target. In the following, Algorithm **MMMu15** has been preferred to Algorithm **MMMu13** because it is more computationally intensive, so the overhead of a correct implementation with respect to one that is linked against an optimized BLAS library is more apparent in the experimental measures. Another reason is that the bound of Theorem 2.1 does not appear directly in Algorithm **MMMu13**.

4 Expressing Parallelism for a Multicore Target

In this section, we present several implementations of **MMMu15** in the C programming language for a multicore target.

The following few assumptions will be made on the structure of the interval matrices in midpoint-radius representation. First, interval matrices are manipulated as pairs of arrays in the row-major order⁴ which is usual in C. Second, the array of midpoints and the array of radii are assumed to share the same alignment and stride. The stride of the pair may be different from the column dimension, as this allows to handle submatrices, but no padding will be made to align rows with a given value.

The C-99 standard defines accesses to the floating-point environment, in particular to rounding modes, through the **fegetround** and **fesetround** functions which allow readings and changes of rounding modes in a portable manner. However, even up-to-date versions of some compilers do not take into account the changes of rounding modes in their optimization phases, yielding possibly incorrect results when they swap a call to **fesetround** and a floating-point operation (cf. the long-running bug of GCC [1]). Calling directly an assembly instruction to change the rounding mode solves this problem, but this workaround is platform-dependent. Nevertheless, in the implementations of **MMMu15** below, we use the portable C-99 functions just before loops, as we did not observe the bug mentioned above in these conditions.

⁴Note that, according to our experience, the **dgemm** function in MKL seems to be more efficient when using the column-major order.

Beside correctness problems, we aim to provide an implementation that is efficient on a multicore target. Two main features of the multiplication of dense matrices make them very well-suited for parallelization. First, they offer a high degree of data independence, as each component of the product could virtually be computed regardless of the others. This offers opportunity for calculating sub-matrix blocks in different threads. Second, the property of density of the matrices permits regularity and contiguity of memory accesses. This last condition is necessary to benefit from the memory cache system of processors.

We exploit below the coarse-grained parallelism of the block calculation at the thread level by using OpenMP constructs. Furthermore, we translate the fine-grained parallelism of contiguous data accesses and similar computations into instruction level parallelism through vector instructions.

Ideally, this last step should be handled by the compiler auto-vectorization, but, as discussed below in Section 4.4, that is beyond current compiler capabilities.

4.1 Version with a Parallel BLAS Library

It has been shown in Section 2 that the algorithms for interval matrix products using the midpoint-radius representation rely on floating-point matrix products. As advocated by their authors, this approach is a straightforward means to benefit from optimized BLAS libraries at a small development cost. From this point of view, **MMu15** can be parallelized on a multicore machine by using a multithreaded BLAS library.

We present below a possible implementation of the first part in rounding-to-nearest mode (lines 1 to 4) of Algorithm 2. In the BLAS implementation of Algorithm 3, lines 1, 4, 7, and 10 translate directly to calls to the **gemm** function, which takes two scalars α and β and three matrices A , B , and C as parameters and computes $C \leftarrow \alpha A \cdot B + \beta C$.

Algorithm 3 Rounding-to-nearest part of **MMu15**-BLAS

Input: $\langle M_A, R_A \rangle$ and $\langle M_B, R_B \rangle$

Output: M_C and Γ

- 1: $M_C \leftarrow 1M_A \cdot M_B + 0M_C$
 - 2: $T_1 \leftarrow |M_A|$
 - 3: $T_2 \leftarrow |M_B|$
 - 4: $\Gamma \leftarrow 1T_1 \cdot T_2 + 0\Gamma$
 - 5: $T_3 \leftarrow \min(T_1, R_A)$
 - 6: $T_4 \leftarrow \min(T_2, R_B)$
 - 7: $\Gamma \leftarrow 1T_3 \cdot T_4 + 1\Gamma$
 - 8: $T_5 \leftarrow \text{sign}(M_A) \cdot T_3$
 - 9: $T_6 \leftarrow \text{sign}(M_B) \cdot T_4$
 - 10: $M_C \leftarrow 1T_5 \cdot T_6 + 1M_C$
-

Alas, as discussed in Section 3.2, lines 1 and 4 compute $M_A \cdot M_B$ and $|M_A| \cdot |M_B|$ with arithmetic operations in arbitrary orders. Therefore, we cannot guarantee that this implementation does always satisfy the Theorem 2.1 hypothesis of identical orders.

4.2 Version with a Parallel Loop for the Rounding-to-Nearest Part

The simplest means of preserving the order of operations between products and products of absolute value is to combine their computations in the same loop. We call **MMul5frn** (for *fused loops in rounding-to-nearest*) the Algorithm 4 below that computes simultaneously M_C and Γ in the classical three nested loops of a matrix product. At lines 8 and 10, note that ef is of the same sign than ab , so $|ab|+|ef| = |ab+ef| = |p|$.

Algorithm 4 Rounding-to-nearest part of **MMul5frn**

Input: $\langle M_A, R_A \rangle$ and $\langle M_B, R_B \rangle$

Output: M_C and Γ accumulated in the *same order*

```

1: for  $i \leftarrow 1, m$  do {OpenMP parallel for}
2:   for  $l \leftarrow 1, k$  do
3:     for  $j \leftarrow 1, n$  do
4:        $a \leftarrow M_{Ail}, c \leftarrow R_{Ail}$ 
5:        $b \leftarrow M_{Blj}, d \leftarrow R_{Blj}$ 
6:        $e \leftarrow \text{sign}(a) \min(|a|, c)$ 
7:        $f \leftarrow \text{sign}(b) \min(|b|, d)$ 
8:        $p \leftarrow ab + ef$ 
9:        $M_{Cij} \leftarrow M_{Cij} + p$ 
10:       $\Gamma_{ij} \leftarrow \Gamma_{ij} + |p|$ 
11:     end for
12:   end for
13: end for

```

The choice of the order of the three loops is important with respect to the locality of the memory accesses. In **MMul5frn**, the inner loop spans the columns of M_C and Γ , so that memory is written contiguously. This facilitates hardware prefetching as well as vectorization. The external loop traverses the rows of M_C and Γ and is annotated with an OpenMP **parallel for** directive to distribute blocks of rows to multiple threads. This choice of indices for the external and inner loops also avoids false sharing⁵ between cache lines of different cores when the matrices are stored in row-major order.

Algorithm **MMul5frn** uses less memory than the BLAS implementation described in Section 4.1, each component of ρ values (lines 1 and 2 of Algorithm 2) and of absolute values of M_A and M_B being computed on-the-fly.

⁵The first levels of memory caches are usually private to the core. When two private caches try to maintain a coherent vision of the same memory region, they use a system of monitoring and copy. Modifying a data in some memory block that is duplicated in another cache triggers this costly mechanism, even if the other core does not use the modified data but a nearby one contained in the same cache line. This phenomenon is called *false sharing* (cf. [5, Section 4.3]).

4.3 Version with a Blocked Parallel Loop for the Rounding-to-Nearest Part

When the matrix C is large enough, a row vector of M_C and a row vector of Γ cannot be kept entirely and simultaneously in the first level cache. As the inner loop accumulates partial products in M_C and Γ , the first components are evicted from the cache when the last ones are computed. At the next iteration, the foremost ones are needed again and in turn evict the end ones. With the previous choice of loop nesting, this leads to one cache miss per component at each iteration over the common dimension of M_A and M_B .

To avoid the cache misses of `MMul5frn`, we propose a blocked implementation `MMul5bfrn` where the matrices M_C and Γ are split into blocks of rows that can be stored simultaneously in the same cache. If the capacity of the cache is too small to contain two complete rows, then the rows of M_C and Γ to be computed are also divided into pieces that are small enough to fit in the cache at the same time.

The `MMul5frn` and `MMul5bfrn` versions only differ in the blocking of m and n dimensions, the outer loop is distributed to threads with the OpenMP `parallel for` construct in both implementations and the inner loop (lines 4 to 10 in Algorithm 4) remains unchanged.

4.4 Version with an Explicitly Vectorized Kernel

The inner loop of Algorithm 4 can be very easily and efficiently translated into a sequence of vector instructions without any conditional branch. In fact, signs of floating-point numbers can be extracted and copied with bit-masks and logical operations, while the minimum of two floating-point values is usually provided either as a vector instruction (e.g. `FMIN` on Sparc), or it can be implemented as a combination of the comparison instruction that issues bit masks of all ones or all zeros and logical operations between these bit masks and the floating-point data (e.g. with `Altivec`, `SSE`, and `AVX` instruction sets).

Unfortunately, the transformations that convert the inner loop into vectorized code are still too complicated to be handled by current compilers. Thus, for the purpose of comparing execution time of optimized computation kernel in Section 5, we manually translated the inner loop block Algorithm `MMul5bfrn` into `AVX` code and we denote by `MMul5bfrn-avx` the new implementation. As expected, the gain in execution time is important (see Section 5.3), but the code is no more portable.

4.5 Rounding-Upwards Part

The computation of R_C uses the rounding to $+\infty$ mode, as shown in the lines 5 and 6 of Algorithm 2. It involves one matrix product and we cannot ensure that an `gemm` function of a given BLAS library always yields the expected overestimation (see Section 3.1). Therefore, if we want to guarantee the inclusion of the exact product in the computed product, an implementation of Algorithm 2 has to be statically linked against a trusted BLAS library if it uses the BLAS `gemm` function. Another solution is to write the rounding-upward part with the classical three loops algorithm.

In the latter case, the implementation can use strategies similar to the ones used above for the rounding-to-nearest part. The computation of lines 5 and 6 of Algorithm 2 can be mixed into the same inner loop, saving some accesses to temporary memory buffers; remaining memory operations can be organized so that memory is

written in a contiguous manner; and the external loop can be easily parallelized by distributing iterations to OpenMP threads.

5 Performance Results

We now present experimental measures and comparisons of the execution time for the implementations presented in the previous section, namely: the non-guaranteed `MMul5-BLAS` (4.1) and the guaranteed versions `MMul5frn` (4.2), `MMul5bfrn` (4.3), and `MMul5bfrn-avx` (4.4). In order to evaluate the costs of the different implementations of the rounding-to-nearest part, all the tested computation kernels share the same rounding-upwards part. This part is not implemented as described in Section 4.5 but simply with a BLAS matrix product (i.e. one call to the `dgemm` function).

5.1 Experimental Setup

We experiment with a shared memory system of 4 eight-cores Xeon E5-4620 processors. To assure stable and reproducible measurement, the Hyper-Threading and Turbo Boost capabilities are disabled and the clock frequency is statically set to its higher possible value. The processor characteristics, operating system, and software being used are described in Table 1.

	clock speed	2.20GHz
	SIMD instruction set	SSE2, AVX
4 processors	level 1 data caches	32KB per core
Intel Xeon E5-4620	level 2 caches	256KB per core
(Sandy Bridge)	level 3 cache	16MB, shared
number of cores	32 (4 processors \times 8 cores)	
operating system	Linux version 3.2.0 (Debian Wheezy)	
compiler suite	GCC 4.7.2	
BLAS library	Intel MKL version 11.0.2	
OpenMP library	Gnu OpenMP library	

Table 1: Description of the system used for performance results.

We compile with GCC version 4.7.2 and the following compiler options: `-O2` sets the level of optimization, `-m64` is required when linking with the 64-bit version of the MKL, `-frounding-math` disables floating-point optimizations that are valid only with the default rounding-to-nearest mode, `-mfpmath=sse` imposes the use of SSE/AVX floating-point unit, `-mavx` enables the AVX-vector instruction set, and `-ftree-vectorize` unconditionally triggers the loop vectorization phase. Figures do not differ noticeably with the `-O3` level of optimization, except for very small matrices.

As it has been discussed in Section 4, the three kernels `MMul5frn`, `MMul5bfrn`, and `MMul5bfrn-avx` are multithreaded using OpenMP. OpenMP libraries are controlled by environment variables that are also taken into account by the MKL library for its own multithreading management. In our experiments, we set the following environment variables: `OMP_PROC_BIND` is set to `true` in order to avoid thread migration from one core to another and the associated time overhead, and `OMP_DYNAMIC` is set to `false` so as

to benefit from all the available cores when wanted without being under the control of dynamic adjustment of the number of threads. Finally, we choose the double precision (i.e. the IEEE-754 **Binary64** type) as the working precision in all measures below.

5.2 Sequential Execution Time

We first evaluate and compare the sequential execution time of the computation kernels presented in Section 4. The measured timings are displayed on Figure 1.

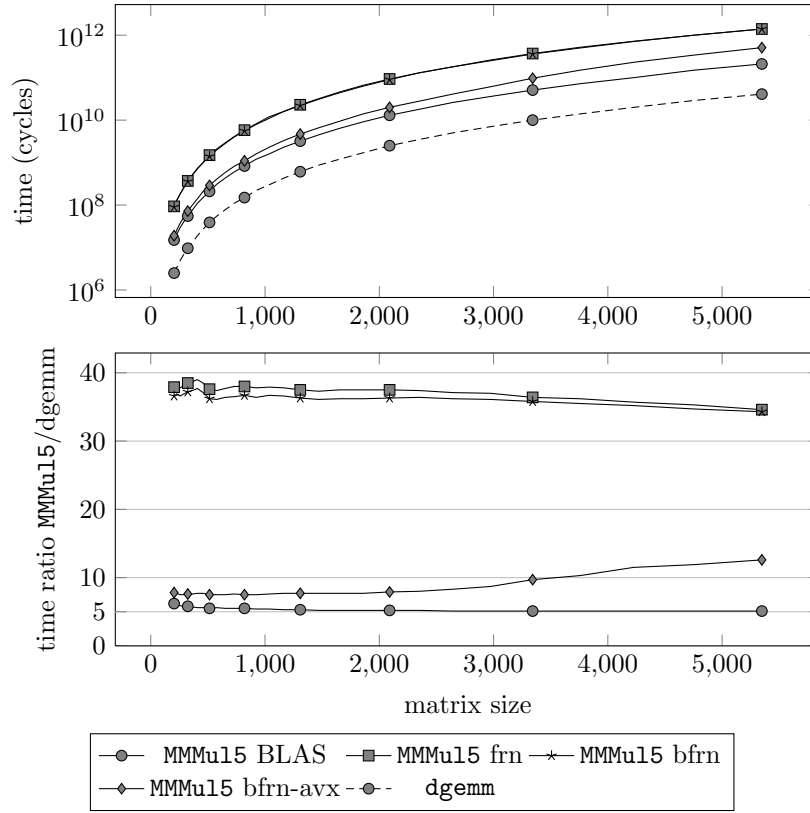


Figure 1: Execution time – Sequential.

The graph at the top of Figure 1 represents the execution time of the interval matrix products along with the floating-point matrix product (**dgemm**). The input is a pair of square matrices from dimension 200 to 5,350. In order to exhibit the slowdown of interval compared to floating-point matrix products, the ratios of execution time of the different implementations to the time of a **dgemm** execution are displayed at the bottom of Figure 1.

The following remarks can be made about these results. As the authors of Algorithm 2 expected, the cost of **MMu15-BLAS** is very close to 5 times the cost of a call to **dgemm**. Others implementations do not benefit from the optimization of the MKL and

the `MMu15frn` is about 35 times slower than `dgemm`. The blocked version `MMu15bfrn` is slightly faster than `MMu15frn` because of a better locality of memory accesses as explained in Section 4.3. The explicit AVX version `MMu15bfrn-avx` is about four times as fast as `MMu15bfrn`, demonstrating that the compiler is unable to efficiently vectorize the code by itself. Indeed, on the `MMu15bfrn` code, we observed little gain provided by the `-mavx` compiler option, because GCC is unable to vectorize the inner loop. For small matrices (up to $3,500 \times 3,500$), `MMu15bfrn-avx` is less than 2 times slower than the non-guaranteed `MMu15-BLAS`.

5.3 Execution Time with 32 Threads

We measure now the execution time of the computation kernels when using one thread per core. The measured timings are displayed on Figure 2 in the placement described in the previous section.

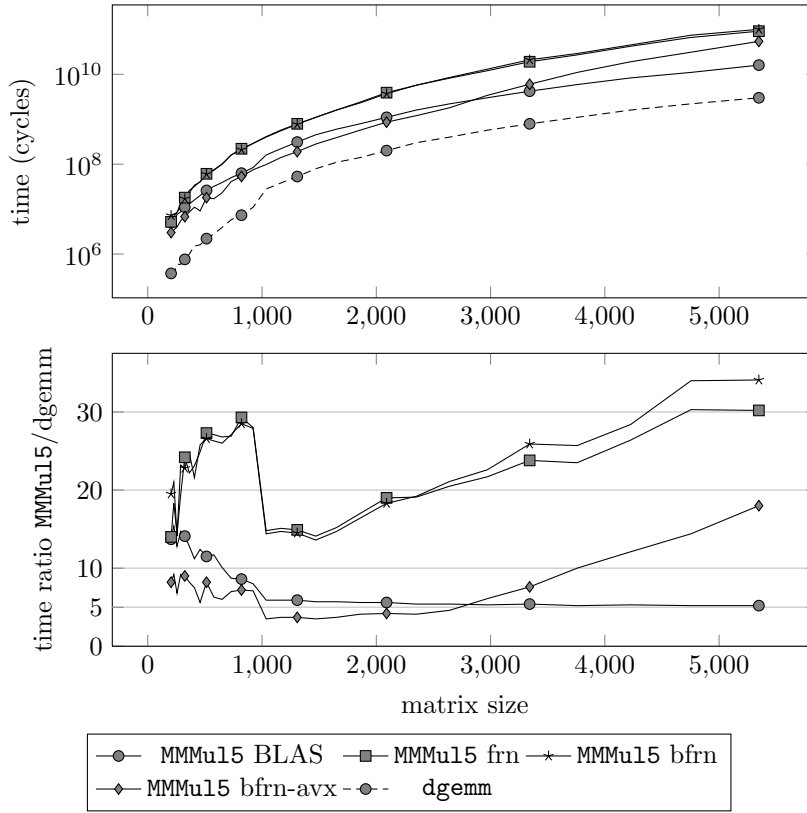


Figure 2: Execution time – 32 threads.

It can be noticed in the graph at the top of Figure 2 that `dgemm` and `MMu15-BLAS` are subject to a sudden increase of execution time around matrix dimension of 1,000. We ascribe this phenomenon to the capacity of the level 3 cache that can contain

simultaneously no more than two square matrices of dimension 1,024. With matrices of dimension larger than 1,000, the slope of the MKL timings tends to be more and more gentle, demonstrating the optimization of the library for large matrices. Kernels `MMu15frn` and `MMu15bfrn` present comparable timings until dimension 2,500. From dimension 500 to dimension 2,000, `MMu15bfrn` undergoes less cache misses when accessing the first level of memory cache and is slightly more efficient than `MMu15frn`. For larger matrices, the blocked version `MMu15bfrn` appears to be the slowest one. This may be due to the greater management overhead of an higher number of OpenMP threads. Finally, this set of measures demonstrates that a blocked implementation correctly translated in vector instructions like `MMu15bfrn-avx` can compete with the BLAS implementation for matrices of dimension less than 2,500.

Compared to the sequential execution, the relative behaviors of the kernels are markedly different, as it can be noticed from the graph at the bottom of Figure 2. For small input, the overhead of non-BLAS operations (absolute values, minimum selections, and sign copies, see Algorithm 3 p.7) slows down `MMu15-BLAS`. This fact is more evident in the heavily parallel execution because the execution time of `dgemm` is one order of magnitude smaller. Then, as the matrix dimension grows, the ratio perceptibly tends to 5. Kernels `MMu15frn` and `MMu15bfrn` are still much less efficient than the BLAS implementation, but to a lesser extend than in the sequential case. The sudden decrease of the ratio value around dimension 1,000 is to be related to the time increase of the reference `dgemm` already discussed.

To compare the best implementations, `MMu15bfrn-avx` is faster than the non-guaranteed `MMu15-BLAS` kernel for matrices up to dimension 2,500. Then, `MMu15bfrn-avx` gets slower and slower: this is due to the fact that the blocking strategy is limited to the first level of cache. An improved version with several levels of blocking and with thresholds tuned to the underlying cache hierarchy should remain competitive with the BLAS implementation for larger matrices.

5.4 Scalability

Finally, we analyze the *strong scalability* of the computation kernels, that is, their behavior when the number of threads rises while the matrix dimension remains constant. One important metric in this context is the measure of *efficiency*. By definition, the efficiency of a parallel implementation executed with p threads is the ratio $T_1/(p \times T_p)$, where T_1 is the sequential execution time, and T_p is the execution time using p execution threads. It equals one when the parallelism is perfect.

In Figure 3 are shown the measures of execution time (left) and scalability (right) of the `MMu15-BLAS` and `MMu15bfrn-avx` kernels, along with the corresponding values for `dgemm` implementation of the MKL. The input factors of the products are square matrices of dimension 1024, and the number of threads varies from 1 to 32.

The analysis of these measures explains why the `MMu15bfrn-avx` implementation, which is slower than `MMu15-BLAS` in sequential execution, turns superior at high level of parallelism. In fact, the former implementation is about 1.4 times slower than the BLAS implementation with one thread, its execution time is comparable with 8 threads, and 1.5 times less with 32 threads. This phenomenon is more explicit in the left part of Figure 3: up to one thread per processor (4 threads), the efficiency of MKL `dgemm` is higher than 90 %, then it falls to low levels, near 30 %, with 32 threads. This is due to the fact that the memory bandwidth is shared among cores in a processor and this resource rapidly becomes a limiting factor as the number of threads augments. The `MMu15-BLAS` kernel, which depends on the MKL `dgemm` for its more intensive

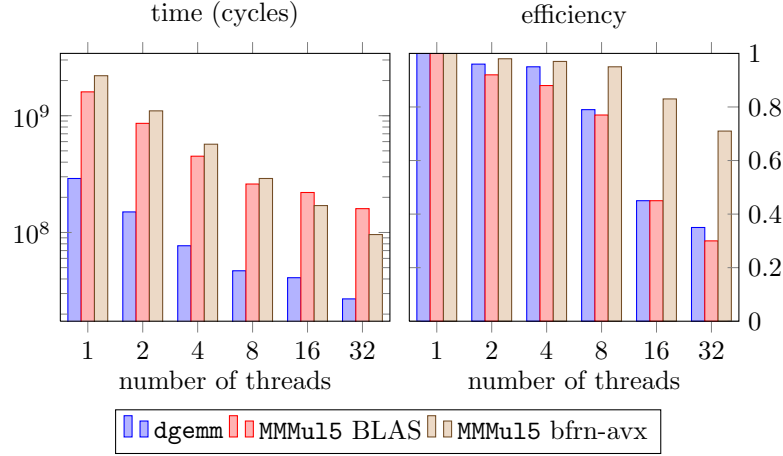


Figure 3: Scalability – 1024-by-1024 matrices.

computational part, exhibits the same important decrease in efficiency. However, the behavior of **MMu15bfrn-avx** demonstrates the possibility of improvement. In contrast to the BLAS implementation, the efficiency of **MMu15bfrn-avx** decreases slowly and never drops below value of about 70 % with 32 threads. The efficiencies of **MMu15frn** and **MMu15bfrn**, not shown here, follow the same scheme. This can be explained by the merging of the computations of $M_{\mathbf{A}} \cdot M_{\mathbf{B}}$, $\rho_{\mathbf{A}} \cdot \rho_{\mathbf{B}}$, $|M_{\mathbf{A}}| \cdot |M_{\mathbf{B}}|$, and $|\rho_{\mathbf{A}}| \cdot |\rho_{\mathbf{B}}|$ in the fused loop versions. It saves 3 successive calls to **dgemm**, which are as many global synchronization points in a parallel execution. On the contrary, the threads in the BLAS implementation have to wait for the latest ones among them, that are slowed down by concurrency to memory resource or system interruptions for instance, before entering the remaining computation.

6 Conclusion

The new results presented in Section 4 are algorithms for interval matrix multiplication that are really guaranteed to contain the exact result and that are efficient on a multicore architecture, as seen in Section 5.

We have seen that implementing interval algorithms on high-performance architectures amplifies the recurrent question of the efficiency of these algorithms, as poor performance even degrades compared to the performance of (non-guaranteed but fast) floating-point algorithms. Indeed, getting reliability and performance is even more difficult than on sequential architectures, as many issues are obscured by lack of specifications (for instance, what happens to the rounding mode in a multithreaded environment) or by specific issues, such as the lack of reproducibility of the numerical result from run to run.

From the implementation developed in this paper, we may draw the following methodological conclusions. Firstly, develop interval algorithms that are based on well-established numerical bricks (such as the **gemm** of BLAS in this work), so as to benefit

from their optimized implementation. A second step could be to convince developers and vendors of these bricks to clearly specify their behavior, especially with regards to rounding mode. However, we have observed that this may not suffice, as this was the case with the issue on the order in which floating-point operations should be performed at each call. The approach adopted in the work presented here has been to replicate the work done for the optimization of the considered numerical bricks, and to adapt it to the specificities and requirements of the interval algorithm, as it has been done here to compute “simultaneously” $A \cdot B$ and $|A| \cdot |B|$. As it can be observed on the experimental results, this approach pays off, as it makes it possible to guarantee the inclusion property – which is the fundamental property and strength of interval algorithms – and to get performances. To obtain even better performances, for the interval matrix multiplication in particular, it would be worth developing an ATLAS-like autotuning of the blocked version in order to optimize the usage of all cache levels and not only of the cache of level 1, as it has been done here. This would prevent the loss of performance when the matrix size increases too much.

However, following such a methodology requires that the programmer gathers skills both in interval analysis and in HPC programming and this is rare.

This methodology could lead to the development of Interval-BLAS, offering various algorithms with different accuracies (at least theoretically established) and good efficiency. The existence of such a library would be valuable to promote the wide use of interval algorithms, while alleviating the burden to get performances, from the user of these algorithms.

References

- [1] Bug 34678 – optimization generates incorrect code with -frounding-math option. http://gcc.gnu.org/bugzilla/show_bug.cgi?id=34678, January 2008.
- [2] S. Blackford, G. Corliss, J. Demmel, J. Dongarra, I. Duff, S. Hammarling, G. Henry, M. Heroux, C. Hu, W. Kahan, L. Kaufmann, B. Kearfott, F. Krogh, X. Li, Z. Maany, A. Petitet, R. Pozo, K. Remington, W. Walster, C. Whaley, J. Wolff von Gudenberg, and A. Lumsdaine. *Basic Linear Algebra Subprograms Technical (BLAST) Forum Standard*, 2001. <http://www.netlib.org/blast-forum/>.
- [3] Kazushige Goto and Robert van de Geijn. High-performance implementation of the level-3 BLAS. *ACM Transactions on Mathematical Software*, 35(1):4:1–4:14, 2008.
- [4] Y. He and C. H. Q. Ding. Using accurate arithmetics to improve numerical reproducibility and stability in parallel applications. *The Journal of Supercomputing*, 18:259–277, 2001.
- [5] John L. Hennessy and David A. Patterson. *Computer Architecture - A Quantitative Approach (4th ed.)*. Morgan Kaufmann, 2007.
- [6] Intel. *Intel Math Kernel Library for Linux OS User's Guide*. http://software.intel.com/sites/products/documentation/doclib/mkl_sa/11/mkl_userguide_lnx/mkl_userguide_lnx.pdf.
- [7] R. B. Kearfott, M. T. Nakao, A. Neumaier, S. M. Rump, S. P. Shary, and P. van Hentenryck. Standardized notation in interval analysis. *Reliable Computing*, 15(1):7–13, 2010.

- [8] Christoph Lauter and Valérie Ménessier-Morain. There's no reliable computing without reliable access to rounding modes. In *SCAN 2012 Symposium on Scientific Computing, Computer Arithmetics and Verified Numerics*, pages 99–100, 2012.
- [9] Xiaoye S. Li, James W. Demmel, David H. Bailey, Greg Henry, Yozo Hida, Jimmy Iskandar, William Kahan, Suh Y. Kang, Anil Kapur, Michael C. Martin, Brandon J. Thompson, Teresa Tung, and Daniel J. Yoo. Design, Implementation and Testing of Extended and Mixed Precision BLAS. *ACM Transactions on Mathematical Software*, 28(2):152–205, June 2002.
- [10] Jean-Michel Muller, Nicolas Brisebarre, Florent de Dinechin, Claude-Pierre Jeanerod, Vincent Lefèvre, Guillaume Melquiond, Nathalie Revol, Damien Stehlé, and Serge Torres. *Handbook of Floating-Point Arithmetic*. Birkhäuser Boston, 2010.
- [11] A. Neumaier. *Interval methods for systems of equations*. Cambridge University Press, 1990.
- [12] Hong Diep Nguyen. *Efficient algorithms for verified scientific computing: numerical linear algebra using interval arithmetic*. PhD thesis, École Normale Supérieure de Lyon - Université de Lyon, 2011. <http://hal-ens-lyon.archives-ouvertes.fr/ensl-00560188>.
- [13] Takeshi Ogita and Shin'ichi Oishi. Fast inclusion of interval matrix multiplication. *Reliable Computing*, 11(3):191–205, 2005.
- [14] Katsuhisa Ozaki, Takeshi Ogita, Siegfried M. Rump, and Shin'ichi Oishi. Fast algorithms for floating-point interval matrix multiplication. *Journal of Computational and Applied Mathematics*, 236:1795–1814, 2012.
- [15] Siegfried M. Rump. Fast and parallel interval arithmetic. *BIT*, 39:534–554, 1999.
- [16] Siegfried M. Rump. Fast interval matrix multiplication. *Numerical Algorithms*, 61(1):1–34, 2012.
- [17] Volker Strassen. Gaussian elimination is not optimal. *Numerische Mathematik*, 13:354–356, 1969.
- [18] R. Todd. Introduction to Conditional Numerical Reproducibility (CNR). <http://software.intel.com/en-us/articles/introduction-to-the-conditional-numerical-reproducibility-cnr>, 2012.
- [19] R. Clint Whaley and Jack Dongarra. Automatically Tuned Linear Algebra Software. In *Conference on High Performance Networking and Computing*, pages 1–27. IEEE Computer Society, 1998.